

## Logistic Regression

\* binary classification and have the probability of each class

1. we need to limit  $w^T x + b$  to  $[0, 1]$ , so use sigmoid function  
↳ probability

$$\hat{y} = \sigma(w^T x + b) \quad \sigma(z) = \frac{1}{1+e^{-z}} \quad \text{for classification: } \hat{y} = \begin{cases} 1 & \text{if } \sigma(z) \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

good for binary classification

2. when you add the sigmoid activation function, you have to make sure that the loss function is convex  
loss

the cross-entropy loss function is common. Why? It is derived from maximum likelihood estimation.

for 1 observation, we are given features  $x_i$ , weights  $\theta$ , and value  $y_i \in \{0, 1\}$

$$\text{probability of } y_i \text{ given } x_i \text{ and } \theta \quad \theta(x_i, \theta) = \frac{1}{1+e^{-(\theta_0 + \theta^T x_i)}} \quad \text{sigmoid}$$

$$P(y=1 | x_i, \theta) = \theta(x_i, \theta) \quad P(y=0 | x_i, \theta) = 1 - \theta(x_i, \theta) \quad ] \text{ recall that there are only 2 classes}$$

$$\downarrow P(y | x_i, \theta) = [\theta(x_i, \theta)]^{y_i} [1 - \theta(x_i, \theta)]^{1-y_i}$$

for  $n$  observations, it is the product of the probability.

we want the weights  $\theta$  that maximize this probability

$$\lambda(x_i, \theta) = \max_{\theta} \prod_{i=1}^n [\theta(x_i, \theta)]^{y_i} [1 - \theta(x_i, \theta)]^{1-y_i}$$

$$\downarrow \hat{y} = \theta(x_i, \theta)$$

$$= \max_{\theta} \prod_{i=1}^n [\hat{y}_i]^{y_i} [1 - \hat{y}_i]^{1-y_i}$$

↓ take log to turn products to sum

$$\text{this is called log loss} \quad \ln(\lambda(x_i, \theta)) = \max_{\theta} \sum_{i=1}^n \ln([\hat{y}_i]^{y_i} [1 - \hat{y}_i]^{1-y_i})$$

$$\ell(\theta) = \max_{\theta} \sum_{i=1}^n y_i \ln(\hat{y}_i) + (1-y_i) \ln(1-\hat{y}_i)$$

$$\downarrow \text{or cross entropy loss function (for binary case)} \quad Y \in \{0, 1\}$$

the  $(y_i)$  and  $(1-y_i)$  terms are like an indicator variable  $\mathbb{1}\{Y_i = k\}$   
 $Y \in \{0, 1\}$

## Multi-class Classification

one versus all

-  $K$  binary classifiers that say it belongs to the class or not, return model with highest probability

all versus all

Softmax Regression essentially a neural network

- extends logistical regression to multiple classes and gives likelihoods of each label

in logistical regression, since there are only 2 classes, we did:  
 $P(Y=1 | x_i, \theta) = a$  and  $P(Y=0 | x_i, \theta) = 1 - a$

Softmax regression when  $K=2$  is logistical regression

- each class out of  $K$  classes has its own weight vector  $w_{ik}$

- score for class  $k$  using weight vector  $w_{ik} : w_{ik}^T x$

probability for class  $k$ :

$$P(Y=k | x, w_{ik}) = \frac{e^{w_{ik}^T x}}{\sum_{j=1}^K e^{w_j^T x}}$$

softmax:  $f(x, w) = \frac{1}{\sum_{j=1}^K e^{w_j^T x}} \begin{bmatrix} e^{w_1^T x} \\ \vdots \\ e^{w_K^T x} \end{bmatrix}$  where  $W = \begin{bmatrix} -w_1^T \\ -w_2^T \\ \vdots \\ -w_K^T \end{bmatrix}$

Sigmoid gradient

$$\nabla_z g(z) = g(z)[1-g(z)]$$

$$g(z_i) = \frac{1}{1+e^{-z_i}}$$

$$\begin{aligned}
\nabla_z g(z) &= \nabla_z \left[ \frac{1}{1+e^z} \right] \\
&= \nabla_z [1+e^z]^{-1} \\
&= -(1+e^z)^{-2} \cdot \nabla_z [1+e^z] \\
&= -(1+e^{-z})^{-2} \cdot (e^{-z}) \cdot \nabla_z (-z) \\
&= (1+e^{-z})^{-2} e^{-z} \\
&\circ \frac{e^{-z}}{(1+e^{-z})^2} \\
&= \frac{1+e^{-z}-1}{(1+e^{-z})^2} \\
&= \frac{1+e^{-z}}{(1+e^{-z})^2} - \frac{1}{(1+e^{-z})^2} \\
&= \frac{1}{1+e^{-z}} - \frac{1}{(1+e^{-z})^2} \\
&= \frac{1}{1+e^{-z}} \left( 1 - \frac{1}{1+e^{-z}} \right) \\
&= g(z_i)(1-g(z_i))
\end{aligned}$$

$$\text{recall } g(z_i) = \frac{1}{1+e^{-z_i}}$$

Logistic Regression Gradient Descent

(lecture 12 notes)

$$J(\theta) = -\sum_{i=1}^n [y_i \ln(h_\theta(x_i)) + (1-y_i) \ln(1-h_\theta(x_i))]$$

We want  $\theta = \theta - \eta \nabla_\theta J(\theta)$  for gradient descent

$$\begin{aligned}
J(\theta) &= -\sum_{i=1}^n [y_i \ln(h_\theta(x_i)) - y_i \ln(1-h_\theta(x_i)) + \ln(1-h_\theta(x_i))] \\
&= -\sum_{i=1}^n [y_i \ln(\frac{h_\theta(x_i)}{1-h_\theta(x_i)}) + \ln(1-h_\theta(x_i))]
\end{aligned}$$

$$\begin{aligned}
\frac{h_\theta(x_i)}{1-h_\theta(x_i)} &\stackrel{\leftarrow}{=} \frac{\frac{1}{1+e^{-(\theta_0+\theta^T x_i)}}}{\frac{e^{-(\theta_0+\theta^T x_i)}}{1+e^{-(\theta_0+\theta^T x_i)}}} \\
&= \frac{1}{\frac{e^{-(\theta_0+\theta^T x_i)}}{1+e^{-(\theta_0+\theta^T x_i)}}} \\
&= \frac{e^{-(\theta_0+\theta^T x_i)}}{e^{-(\theta_0+\theta^T x_i)} + 1}
\end{aligned}$$

$$h_\theta(x_i) = \frac{1}{1+e^{-(\theta_0+\theta^T x_i)}}$$

$$\begin{aligned}
J(\theta) &= -\sum_{i=1}^n [y_i \ln(e^{\theta_0+\theta^T x_i}) + \ln(\frac{1}{1+e^{\theta_0+\theta^T x_i}})] \\
J(\theta) &= -\sum_{i=1}^n [y_i (\theta_0 + \theta^T x_i) - \ln(1+e^{\theta_0+\theta^T x_i})]
\end{aligned}$$

$$\nabla_\theta J(\theta) = -\sum_{i=1}^n (\hat{y}_i - y_i) x_{ij}$$

$\hat{y} = h_\theta(x_i)$  in terms of weight vector  $\theta$  and the input  $x_i$

$$\text{recall } \nabla_\theta J(\theta) = \begin{bmatrix} \frac{\partial}{\partial \theta_1} [J(\theta)] \\ \frac{\partial}{\partial \theta_2} [J(\theta)] \\ \vdots \\ \frac{\partial}{\partial \theta_n} [J(\theta)] \end{bmatrix}$$

$$\begin{aligned}
\nabla_\theta J(\theta) &= -\sum_{i=1}^n [y_i \cdot \nabla_\theta (\theta_0 + \theta^T x_i) - \nabla_\theta \ln(1+e^{\theta_0+\theta^T x_i})] \\
&= -\sum_{i=1}^n [y_i \cdot x_i - (\frac{1}{1+e^{\theta_0+\theta^T x_i}}) \cdot \nabla_\theta (e^{\theta_0+\theta^T x_i} + 1)] \\
&= -\sum_{i=1}^n [y_i \cdot x_i - (\frac{1}{1+e^{\theta_0+\theta^T x_i}}) \cdot (e^{\theta_0+\theta^T x_i}) \cdot \nabla_\theta (\theta_0 + \theta^T x_i)] \\
&= -\sum_{i=1}^n [y_i \cdot x_i - x_i \cdot (\frac{e^{\theta_0+\theta^T x_i}}{1+e^{\theta_0+\theta^T x_i}})] \\
&= -\sum_{i=1}^n [y_i \cdot x_i - x_i \cdot (\frac{1}{1+e^{-(\theta_0+\theta^T x_i)}})] \\
&= -\sum_{i=1}^n [y_i \cdot x_i - x_i \cdot h_\theta(x_i)]
\end{aligned}$$

$$\hat{y} = h_\theta(x_i) = \frac{1}{1+e^{-(\theta_0+\theta^T x_i)}}$$

$$= -\sum_{i=1}^n [y_i \cdot x_i - x_i \cdot \hat{y}_i]$$

$$\nabla_\theta J(\theta) = \sum_{i=1}^n x_i (\hat{y}_i - y_i)$$

multiply by  $\frac{1}{n}$  to average out the loss

$$\theta_j = \theta_j - \eta (\sum_{i=1}^n [\hat{y}_i - y_i] x_{ij})$$

$$\begin{aligned}
\theta_0 + \theta^T x_i &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \\
\theta &= \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \quad x_i = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}
\end{aligned}$$

$$\begin{aligned}
\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n &= \theta_0 + \theta_1 x_1 + \dots + \theta_n x_0 \\
\nabla_\theta [\theta_0 + \theta^T x_i] &= x_i
\end{aligned}$$